

# Censored linear model in high dimensions

*Penalised linear regression on high-dimensional data with left censored response variable*

Patric Müller and Sara van de Geer  
Seminar für Statistik, ETH Zurich

## Abstract

Censored data are quite common in statistics and have been studied in depth in the last years (for some early references, see Powell (1984), Murphy et al. (1999), Chay and Powell (2001)). In this paper we consider censored high-dimensional data.

High-dimensional models are in some way more complex than their low-dimensional versions, therefore some different techniques are required. For the linear case appropriate estimators based on penalised regression, have been developed in the last years (see for example Bickel et al. (2009), Koltchinskii (2009)). In particular in sparse contexts the  $l_1$ -penalised regression (also known as LASSO) (see Tibshirani (1996), Bühlmann and van de Geer (2011) and reference therein) performs very well. Only few theoretical work was done in order to analyse censored linear models in a high-dimensional context.

We therefore consider a high-dimensional censored linear model, where the response variable is left-censored. We propose a new estimator, which aims to work with high-dimensional linear censored data. Theoretical non-asymptotic oracle inequalities are derived.

# 1 Introduction

Censored data are quite common in statistics and have been studied in depth in the last years (for some early references, see Powell (1984), Murphy et al. (1999)). We consider the censored linear model, where the response variable is left censored and its non-censored version linearly depends on the predictors. Observed are the predictors, the censored response variable and the censoring level. Our main goal is to recover the linear dependency between the covariates and the uncensored response variable. An example fitting this model is the Social Security Administration earnings records in the years '60s are censored at the 'taxable maximum', that is anyone earning more than the maximum is recorded as having earned at the maximum (see Chay and Powell (2001)). In particular, Powell (1984) proposed an estimator and proved its strong consistency.

In a high-dimensional context, where the dimension of the parameter set  $p$  is bigger than the number of observations  $n \ll p$  the method of Powell is not directly applicable because it would be underdefined. In order to overcome the difficulties given by the high-dimensional case a new estimator is required.

There is a large body of work on linear high-dimensional models. A common approach is to construct penalised estimators like the LASSO (least absolute shrinkage and selection operator proposed in Tibshirani (1996)) or the ridge regression and elastic net (see Zou and Hastie (2005)). The LASSO is widely studied (see e.g. Koltchinskii (2011) and Bühlmann and van de Geer (2011) and references therein) and gives remarkable results in sparse contexts.

In this paper LASSO techniques are combined with the ideas of Powell (1984) in order to obtain a pertinent estimator for the high-dimensional censored linear model. We prove theoretical results and give oracle bounds for both the prediction and the estimation error of our estimator. Simulation supporting the theoretical results are also presented.

The paper is organised as follows. We begin in Section 2 with the description of the model and the required notation and assumptions for the main theorem, which is given in Section 3. The proof of the theorem and the required technical tools can be found in Section 4. Finally in Section 5 we present simulations.

## 2 Model description and notation

Let  $(x, c) \in \mathbb{R}^p \times \mathbb{R}$  be a regression random vector with distribution  $Q_{x,c}$  and  $\varepsilon \in \mathbb{R}$  be an error term with cumulative distribution function  $\nu_0$ . Moreover,  $(x_i, c_i)$  and  $\varepsilon_i$ ,  $i = 1, \dots, n$  are i.i.d. copies of  $(x, c)$  and  $\varepsilon$  respectively.

Consider the following left-censored linear model:

$$y_i = \max \{c_i, x_i \beta^0 + \varepsilon_i\} \quad i = 1, \dots, n. \quad (1)$$

The dependent variable  $y_i$ , the regression vector  $x_i$  and the censoring level  $c_i$  are observed for each  $i$ , while the (conformable) parameter vector  $\beta^0 \in \mathbb{R}^p$  and the error term  $\varepsilon_i$  are unobserved.

**Remark 1.** *Real datasets often have a fixed, known censoring level (e.g. the tax example of Chay and Powell (2001)).*

The following special cases of our model are of special interest:

The **constant-censored model** corresponds to the special case of Model (1), where  $c_i \equiv c_0$  constant and known.

$$y_i = \max \{c_0, x_i \beta^0 + \varepsilon_i\} \quad , \quad i = 1, \dots, n. \quad (2)$$

Furthermore we define the **zero-censored model** as

$$y_i = \max \{0, x_i \beta^0 + \varepsilon_i\} \quad , \quad i = 1, \dots, n. \quad (3)$$

I.e.  $c_i$  in model (1) is fixed and equal 0.

In the low-dimensional case ( $p \ll n$ ) Powell (1984) showed that in Model (3), under some standard assumptions, the estimator

$$\hat{\beta}^{Powell} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n |y_i - \max \{0, x_i \beta\}| \quad (4)$$

is strongly consistent.

Consider now the high-dimensional left-censored linear model, i.e. Model (1) in the high-dimensional case where  $p \gg n$ . Estimator (4) leads now to an underspecified system and can not be used in high-dimensional contexts. A new estimator is therefore required. A common approach to high-dimensional linear data is the so-called  $l_1$ -penalised regression (LASSO). Combining the  $l_1$ -penalty with the idea of Powell (1984) we define a new estimator for  $\beta^0$  in the high-dimensional context.

$$\hat{\beta} := \arg \min_{\beta \in \mathcal{B}} \left\{ \frac{1}{n} \sum_{i=1}^n |y_i - \max \{c_i, x_i \beta\}| + \lambda \|\beta\|_1 \right\} \quad , \quad (5)$$

where  $\mathcal{B}$  is a bounded set.

The idea behind this estimator is that the first term controls the prediction error, whether the second term keeps the sparsity under control. The parameter  $\lambda$  is a trade-off parameter. An appropriate choice of  $\lambda$  results from Theorem 3.2.

## Notation

Hereafter we list some important notation we use in this paper.

We denote by  $x_{ij}$  the  $j$ -th component of the predictor vector  $x_i$ .

Define now for some real function  $f$  on  $\mathbb{R}^p \times \mathbb{R}$  the loss function  $\rho_f$  and the theoretical risk  $P$  as:

$$\rho_f(x, y, c) := |y - f(x, c)|,$$

$$P_{\rho_f} := \mathbf{E} [\rho_f(x, y, c)].$$

The empirical risk is then

$$P_{n, \rho_f} := \frac{1}{n} \sum_{i=1}^n \rho_f(x_i, y_i, c_i).$$

Furthermore define

$$f_0(x, c) := \arg \min_a \mathbf{E} [|y - a| | x, c] \quad (6)$$

and for some  $\beta \in \mathbb{R}^p$

$$f_\beta(x, c) := x^T \beta \vee c.$$

The excess risk for  $f_\beta$  is then

$$\mathcal{E}(f) := P_{\rho_f} - P_{\rho_{f_0}}.$$

Denote with  $\|\beta\|_1 := \sum_{j=1}^p |\beta_j|$  the  $l_1$ -norm of the vector  $\beta$  and with  $S_\beta := \{j : \beta_j \neq 0\}$  its active set, which has cardinality  $s_\beta$ . Then  $S_\beta^c = \{j : \beta_j = 0\}$ . Let then  $\beta_{j,S} := \beta_j \cdot \mathbb{1}_{j \in S}$ ,  $j = 1, \dots, p$ . Finally define for some function  $f$ ,  $\|f\|^2 := \mathbf{E} [f_\beta^2(x, c)]$ .

## Model Assumptions

We now list some conditions we require in order to prove our results.

**Condition 2.1** (Design condition).

For some constant  $K_X$  it holds that

$$\max_{i,j} |x_{ij}| \leq K_X.$$

A bound on the  $X$ -values is a somewhat restrictive assumption. However we can often approximate an unbounded distribution with its truncated version. Furthermore this condition is quite standard in high-dimensional contexts (see Bühlmann and van de Geer (2011)).

**Condition 2.2** (Design condition II).

For some constant  $K_0$  it holds that  $x_i(\beta - \beta^0)$  and  $f_0(x_i) - f_\beta(x_i)$  takes values in interval  $[-K_0, K_0]$  for all  $i = 1, \dots, n, \forall \beta \in \mathcal{B}$ .

**Condition 2.3** (Scaling condition).

$$\mathbf{E} [x_{ij}^2] = 1 \quad \text{for all } j = 1 \dots p, \ i = 1, \dots, p.$$

This condition can be obtained by rescaling the variables. (Note that we assumed that  $x_i$  are i.i.d. for  $i = 1, \dots, p$ ).

**Condition 2.4** (Solution uniqueness).

The function  $f_0(x, c)$ , defined in (6), is uniquely defined.

**Condition 2.5** (Error assumptions).

The distribution function  $\nu_0$  of the error term  $\varepsilon$ , has median 0 and is everywhere continuously differentiable, with Lipschitz derivative  $\dot{\nu}_0$ . Furthermore assume that  $\dot{\nu}_0(0) > 0$

**Condition 2.6** (Censoring condition).

There exists some constant  $C_2 > 0$  such that for all  $\beta$  satisfying  $\|(\beta - \beta^0)_{S_{\beta^0}^c}\|_1 \leq 3\|(\beta - \beta^0)_{S_{\beta^0}}\|_1$  it holds that:

$$\|f_\beta - f_0\|_2^2 \geq C_2 \|x^T(\beta - \beta^0)\|_2^2.$$

This condition only guaranties that there are not, too many censored data. If  $(x_i, c_i)$  are i.i.d. and have a symmetric joint distribution then the censoring condition holds for any  $\beta, \beta^0$  with constant  $C_2 = 1/4$ .

**Condition 2.7** (Compatibility condition).

The **Compatibility condition** is satisfied for the set  $S_{\beta^0}$  if for some  $\phi_0 > 0$  and all  $\beta$  satisfying  $\|(\beta - \beta^0)_{S_{\beta^0}^c}\|_1 \leq 3\|(\beta - \beta^0)_{S_{\beta^0}}\|_1$  it holds that:

$$\|(\beta - \beta^0)_{S_{\beta^0}}\|_1^2 \leq (\beta - \beta^0)^T \mathbf{E} [x^T x] (\beta - \beta^0) \frac{s_{\beta^0}}{\phi_0^2}.$$

Here  $\phi_0^2$  is the so called compatibility constant introduced by van de Geer (2007).

### 3 Results

Consider the high-dimensional version of model (1) ( $p \gg n$ ). Recall the definition of the  $l_1$ -penalised estimator (5).

$$\hat{\beta} = \arg \min_{\beta \in \mathcal{B}} \left[ \frac{1}{n} \sum_{i=1}^n |y_i - \max \{c_i, x_i \beta\}| + \lambda \|\beta\|_1 \right],$$

**Theorem 3.1.** Assume Conditions 2.1-2.7 and define

$$\lambda(t) := 4K_X \sqrt{\frac{2\log(2p)}{n}} + K_X \sqrt{\frac{8t}{n}}.$$

Then for  $\lambda \geq 4\lambda(t)$  and some constant  $C$  depending only on  $K_0$  and  $C_2$  (see Conditions 2.2 and 2.7 respectively), with at least probability  $1 - 1/p$ , it holds

$$\mathcal{E}(f_{\hat{\beta}}) \leq \lambda^2 \frac{9s_{\beta^0} C}{\phi_0^2} \quad (7)$$

and

$$\|\hat{\beta} - \beta^0\|_1 \leq \lambda \frac{6Cs_{\beta^0}}{\phi_0^2}. \quad (8)$$

**Remark 2** (Asymptotics).

Asymptotically we have

$$\mathcal{E}(f_{\hat{\beta}}) = O\left(\frac{s_{\beta^0} \log p}{n}\right)$$

and

$$\|\hat{\beta} - \beta^0\|_1 = O\left(s_{\beta^0} \sqrt{\frac{\log p}{n}}\right)$$

In a sparse context, where

$$\frac{s_{\beta^0} \log p}{n} \rightarrow 0$$

as  $n$ ,  $p$ , and possibly also  $s_{\beta^0}$  tend to infinity, the excess risk converges to 0.

If furthermore  $s_{\beta^0} \sqrt{\frac{\log p}{n}} \rightarrow 0$ , then also the estimation error converges to 0.

**Remark 3.** The bounds for prediction and estimation error given in Theorem 3.1 do not depend on the distribution of the censoring factor. In fact the censoring level  $c$  has a direct influence on the constant  $C_2$  in Assumption 2.6. In general higher values for  $c_i$  increase the number of censored data. This leads to smaller  $C_2$ . The fact that the censoring level does not directly appear in the theorem should be understood in the sense that the percentage of censored data is important, not the censoring level.

We now shortly focus on Model (3), i.e. the special case where  $c_i$  are all fixed and equal 0. In this case the function  $f$ , as well as the loss function do not any more depend on  $c$ . Therefore we just write  $f(x)$  and  $\rho_f(x, y)$  in spite of  $f(x, 0)$  and  $\rho_f(x, y, 0)$  respectively.

The estimator (5) can be rewritten as

$$\hat{\beta}^{res} := \arg \min_{\beta \in \mathcal{B}} \left\{ \frac{1}{n} \sum_{i=1}^n |y_i - \max\{0, x_i \beta\}| + \lambda \|\beta\|_1 \right\},$$

This corresponds to the high-dimensional penalised version of the estimator proposed by Powell.

**Corollary 3.2.** *Assume the same conditions and use the same definitions as in Theorem 3.1, then for  $\lambda \geq 4\lambda(t)$  and some constant  $C$ , with at least probability  $1 - 1/p$ , it holds*

$$\mathcal{E}(f_{\hat{\beta}}) \leq \lambda^2 \frac{9s_{\beta^0}C}{\phi_0^2} \quad (9)$$

and

$$\|\hat{\beta}^{res} - \beta^0\|_1 \leq \lambda \frac{6Cs_{\beta^0}}{\phi_0^2}. \quad (10)$$

The proof directly follows from Theorem 3.1 taking  $c_i \equiv 0$ .

## 4 Proofs

### Preliminary remarks and lemmas

Denote by  $\nu(y|x, c)$  the distribution function of the censored random variable  $y$  given  $x$  and  $c$ . Then

$$\nu(y|x, c) = \begin{cases} 0 & \text{if } y < c \\ \nu_0(y - x\beta^0) & \text{if } y \geq c \end{cases}.$$

Consequently  $\nu(y|x, c)$  is everywhere differentiable, up to  $y = c$ .

Furthermore we show that  $f_0(x, c) = f_{\beta^0}(x, c)$ .

*Proof.*  $\mathbf{E} [|y - a| | x, c]$  is minimized by  $a = \nu^{-1}(\frac{1}{2} | x, c)$  (the median of  $y$  given  $x$  and  $c$ ). We have:

$$\begin{aligned} f_0(x, c) &= \\ \text{median}(y|x, c) &= \text{median}(\max\{x\beta^0 + \varepsilon, c\} | x, c) \\ &= \max\{x\beta^0 + \text{median}(\varepsilon), c\} \\ &= \max\{x\beta^0, c\} \quad (= x\beta^0 \vee c) \\ &= f_{\beta^0}(x, c). \end{aligned}$$

□

**Remark 4.** *Because  $f_0$  is a minimizer of the excess risk,*

$$\mathcal{E}(f) \geq 0 \quad \forall f.$$

**Lemma 4.1.** *Assume Conditions 2.4 and 2.5, then for all  $\beta \in \mathcal{B}$*

$$\mathcal{E}(f_\beta) \geq C_1^2 \|f_\beta - f_0\|_2^2.$$

*Proof.* For any  $a \geq c$  we have:

$$\begin{aligned} \mathbf{E} \left[ |y - a| \middle| x, c \right] &= \mathbf{E} \left[ (y - a) \mathbb{1}\{y > a\} \middle| x, c \right] - \mathbf{E} \left[ (y - a) \mathbb{1}\{y \leq a\} \middle| x, c \right] \\ &= \mathbf{E} [y - a | x, c] - 2 \mathbf{E} \left[ (y - a) \mathbb{1}\{y \leq a\} \middle| x, c \right] \\ &= \mathbf{E} [y | x, c] - a + 2a\nu(a | x, c) - 2 \mathbf{E} \left[ y \mathbb{1}\{y \leq a\} \middle| x \right]. \end{aligned}$$

Thus

$$\begin{aligned} &\mathbf{E} \left[ |y - a| \middle| x \right] - \mathbf{E} \left[ |y - f_0(x, c)| \middle| x \right] \\ &= f_0 - a + 2a\nu_0(a - x\beta^0) - 2f_0\nu_0(f_0 - x\beta^0) \\ &\quad - 2 \int_{f_0 - x\beta^0}^{a - x\beta^0} (x\beta^0 + \varepsilon) d\nu_0(\varepsilon). \end{aligned} \tag{11}$$

Define  $z := a - f_0$  and first look at the case  $f_0 > c$ , i.e. we have  $f_0 = x\beta^0$ . Then the above expression can be rewritten as:

$$\begin{aligned} (11) &= -z + 2z\nu_0(z) - 2 \int_0^z \varepsilon d\nu_0(\varepsilon) \\ &= 2 \int_0^z (z - \varepsilon) d\nu_0(\varepsilon) \\ &= 2 \int_0^z (z - \varepsilon) \dot{\nu}_0(0) d\varepsilon + 2 \int_0^z (z - \varepsilon) (\dot{\nu}_0(\varepsilon) - \dot{\nu}_0(0)) d\varepsilon. \end{aligned} \tag{12}$$

Using the lipshitz condition in the second integral and integrating we finally obtain:

$$(12) \geq z^2 \dot{\nu}_0(0) - \frac{L}{3} |z|^3. \tag{13}$$

If  $x\beta < c$ , then  $f_0(x, c) = c$ . Define  $q := c - x\beta$ . Expression (11) can be



simplified as follows:

$$\begin{aligned}
& -z + 2a\nu_0(z+q) - 2c\nu_0(q) - 2 \int_q^{z+q} (x\beta^0 + \varepsilon) d\nu_0(\varepsilon) \\
&= 2 \left[ \int_0^{z+q} z d\nu_0(\varepsilon) + \int_q^{z+q} q d\nu_0(\varepsilon) - \int_q^{z+q} \varepsilon d\nu_0(\varepsilon) \right] \\
&= 2 \left[ \int_0^z (z - \varepsilon) d\nu_0(\varepsilon) \right. \\
&\quad \left. + \int_0^{z \wedge q} \varepsilon d\nu_0(\varepsilon) + \int_{z \wedge q}^{z \vee q} z \wedge q d\nu_0(\varepsilon) + \int_{z \wedge q}^{z+q} (z + q - \varepsilon) d\nu_0(\varepsilon) \right] \\
&\geq 2 \int_0^z (z - \varepsilon) d\nu_0(\varepsilon) + 0 + 0 + 0 \\
&\geq z^2 \dot{\nu}_0(0) - L \frac{L}{3} |z|^3,
\end{aligned}$$

where in the last two steps we used that  $q, z \geq 0$  and result (12). Resuming, for any  $f_0, a \geq c$  we have

$$\mathbf{E} \left[ |y - a| \middle| x, c \right] - \mathbf{E} \left[ |y - f_0(x, c)| \middle| x, c \right] \geq (a - f_0)^2 \dot{\nu}_0(0) - \frac{L}{3} |a - f_0|^3.$$

Define now  $h_{x,c}(z) := \mathbf{E} \left[ |y - a_0 + z| \middle| x, c \right] - \mathbf{E} \left[ |y - a_0| \middle| x, c \right]$ ,  $\Lambda^2 := \dot{\nu}_0(0)$  and  $C := \frac{L}{3}$ . Because of the supposed uniqueness of the minimum we have that

$$\forall \epsilon > 0 \quad \exists \alpha_\epsilon > 0 \text{ such that } \inf_{\epsilon \leq |z| \leq K_0} h_{x,c}(z) > \alpha_\epsilon.$$

Furthermore

$$\forall |z| \leq K_0, \quad h_{x,c}(z) \geq \Lambda^2 z^2 - C |z|^3.$$

The function  $h_{x,c}(z)$  satisfies all assumptions of Lemma 4.2 (see after). Then applying the lemma we obtain:

$$\begin{aligned}
h_{x,c}(z) &\geq C_1^2 z^2 && \text{or equivalently,} \\
\mathbf{E} \left[ |y - a| \middle| x, c \right] - \mathbf{E} \left[ |y - f_0(x)| \middle| x, c \right] &\geq C_1^2 (a - f_0)^2,
\end{aligned}$$

where  $C_1^2$  is defined in the cited lemma. Remark that  $C_1$  does not depend on  $(x, c)$ .

Using the iterated expectation we obtain

$$\begin{aligned}
& \mathbf{E} \left[ \mathbf{E} \left[ |y - f_\beta(x, c)| \middle| x, c \right] - \mathbf{E} \left[ |y - f_0(x, c)| \middle| x, c \right] \right] \\
& \geq C_1^2 \mathbf{E} \left[ (f_\beta(x, c) - f_0(x, c))^2 \right] \\
& \Leftrightarrow P \rho_{f_\beta} - P \rho_{f_0} \geq C_1^2 \|f_\beta - f_0\|_2^2.
\end{aligned}$$

This concludes the proof of Lemma 4.1.  $\square$

**Auxiliary Lemma 4.2** (Städler et al. (2010)). *Let  $h : [-K_0, K_0] \rightarrow [0, \infty[$  have the following properties:*

- $\forall \epsilon > 0 \exists \alpha_\epsilon > 0$  such that  $\inf_{\epsilon < |z| \leq K_0} h(z) > \alpha_\epsilon$ ,
- $\exists \Lambda > 0, C > 0$ , such that  $\forall |z| \leq K_0, h(z) \geq \Lambda^2 z^2 - C|z|^3$ .

Then  $\forall |z| \leq K_0$

$$h(z) \geq C_1^2 z^2$$

where

$$C_1^2 := \min \left\{ \epsilon_0; \frac{\alpha_{\epsilon_0}}{K_0^2} \right\}, \quad \epsilon_0 = \frac{\Lambda^2}{2C}.$$

**Lemma 4.3** (Concentration inequality).

Define

$$\gamma(y, c, x) := |y - x\beta \vee c| - |y - x\beta^0 \vee c|,$$

$$Z_M := \sup_{\|\beta - \beta^0\|_1 \leq M} \left| \frac{1}{n} \sum_{i=1}^n \gamma(y_i, c_i, x_i) - \mathbf{E}[\gamma(y_i, c_i, x_i)] \right|$$

then we have

$$P[Z_M \geq M\lambda(t)] \leq \exp(-t).$$

*Proof.* By Massart's inequality (Theorem 14.2 in Bühlmann and van de Geer (2011)) we have, for any  $t > 0$ :

$$P \left[ Z_M \geq \mathbf{E}[Z_M] + MK_X \sqrt{\frac{8t}{n}} \right] \leq \exp(-t).$$

By Lemma 14.20 in Bühlmann and van de Geer (2011) (contraction inequality) we have:

$$\mathbf{E}[Z_M] \leq 4M \sqrt{\frac{2 \log(2p)}{n}} \cdot K_X.$$

Consequently, for all  $t > 0$  and  $M > 0$

$$P \left[ Z_M \geq 4MK_X \sqrt{\frac{2 \log(2p)}{n}} + MK_X \sqrt{\frac{8t}{n}} \right] \leq \exp(-t)$$

or

$$P[Z_M \geq M\lambda(t)] \leq \exp(-t).$$

□

**Lemma 4.4** (Peeling device).

Define for some  $\delta \geq 0$

$$Z_M^\delta := \sup_{\|\beta - \beta^0\|_1 \leq M} \frac{|\nu_n(\beta) - \nu_n(\beta^0)|}{\|\beta - \beta^0\|_1 \vee \delta},$$

then

$$\Pr(Z_M^\delta > 2\lambda(t)) \leq \log_2 \left( \frac{\lceil \log_2 M \rceil}{\lfloor \log \delta \rfloor} \right) e^{-t}.$$

*Proof.* For the proof we use a Peeling device argument (see Bühlmann and van de Geer (2011) and van de Geer (2000)).

$$\begin{aligned} P(Z_M^\delta > 2\lambda(t)) &= P \left( \sup_{\|\beta - \beta^0\|_1 \leq M} \frac{|\nu_n(\beta) - \nu_n(\beta^0)|}{\|\beta - \beta^0\|_1 \vee \delta} > 2\lambda(t) \right) \\ &\leq \sum_{\lfloor j = -\log_2 M \rfloor}^{\lceil -\log \delta - 1 \rceil} P \left( \sup_{2^{-j-1} \leq \|\beta - \beta^0\|_1 \leq 2^{-j}} \frac{|\nu_n(\beta) - \nu_n(\beta^0)|}{\|\beta - \beta^0\|_1} > 2\lambda(t) \right) \\ &\quad + P \left( \sup_{\|\beta - \beta^0\|_1 \leq \delta} \frac{|\nu_n(\beta) - \nu_n(\beta^0)|}{\delta} > 2\lambda(t) \right) \\ &\leq \sum_{\lfloor j = -\log_2 M \rfloor}^{\lceil -\log \delta - 1 \rceil} P \left( \sup_{2^{-j-1} \leq \|\beta - \beta^0\|_1 \leq 2^{-j}} |\nu_n(\beta) - \nu_n(\beta^0)| > 2^{-j} \lambda(t) \right) + e^{-t} \\ &\leq \sum_{\lfloor j = -\log_2 M \rfloor}^{\lceil -\log \delta - 1 \rceil} e^{-t} + e^{-t} \\ &= (\lceil \log_2 M \rceil - \lfloor \log \delta \rfloor) e^{-t}. \end{aligned}$$

□

*Proof of Theorem 3.1.*

We first show inequality (9).

$$\begin{aligned} P_{\rho_{f_{\hat{\beta}}}} - P_{\rho_{f_0}} &= - \left( P_n - P \right) \left( \rho(f_{\hat{\beta}}) - \rho(f_{\beta^0}) \right) \\ &\quad + P_n(\rho(f_{\hat{\beta}})) + \lambda \|\hat{\beta}\|_1 - P_n(\rho(f_{\beta^0})) - \lambda \|\beta^0\|_1 + \\ &\quad + \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}\|_1 \end{aligned} \tag{14}$$

We now look separately at any line of the right part of the last equation.

By definition of  $\hat{\beta}$ , the evaluation of the second line is never positive. Recall,

$$\begin{aligned} S(\beta) &= \{j | \beta_j \neq 0\} , \\ S^c(\beta) &= \{j | \beta_j = 0\} \text{ and} \\ s_\beta &= \#S(\beta) , \end{aligned}$$

then

$$\|\hat{\beta} - \beta^0\|_1 = \sum_{j \in S(\beta^0)} |\hat{\beta}_j - \beta_j^0| + \sum_{j \in S^c(\beta^0)} |\hat{\beta}_j|.$$

Thus

$$\begin{aligned} \lambda \|\beta^0\|_1 - \lambda \|\hat{\beta}\|_1 &= \lambda \left( \sum_{j \in S(\beta^0)} |\beta_j^0| - \sum_{j \in S(\beta^0)} |\hat{\beta}_j| - \sum_{j \in S^c(\beta^0)} |\hat{\beta}_j| \right) \\ &\leq \lambda \sum_{j \in S(\beta^0)} |\hat{\beta} - \beta^0| - \lambda \sum_{j \in S^c(\beta^0)} |\hat{\beta}_j|. \end{aligned}$$

It is easy to show that  $\|\hat{\beta} - \beta^0\|_1 \ll n$ . Lemma 4.4 leads to

$$\left| - \left( P_n - P \right) \left( \rho(f_{\hat{\beta}}) - \rho(f_{\beta^0}) \right) \right| \leq 2\lambda(t) \cdot (\|\hat{\beta} - \beta^0\|_1 \vee p^{-2}),$$

where the above inequality holds with probability at least  $1 - 1/p$  and is obtained by choosing  $t := 2 \log p$  and  $\delta := p^{-2}$  in the above cited lemma.

If  $\|\hat{\beta} - \beta^0\|_1 \leq p^{-2}$  the estimator is very close to the true parameter. Inequalities (9) and (10) trivially follows from  $p^{-2} \ll \lambda/\phi_0^2$ .

If  $\|\hat{\beta} - \beta^0\|_1 \geq p^{-2}$  then

$$\begin{aligned} \mathcal{E}(f_{\hat{\beta}}) &\leq (2\lambda(t) + \lambda) \cdot \sum_{j \in S(\beta^0)} |\hat{\beta}_j - \beta_j^0| + (2\lambda(t) - \lambda) \sum_{j \in S^c(\beta^0)} |\hat{\beta}_j| \quad (15) \\ &\leq (2\lambda(t) + \lambda) \cdot \sum_{j \in S(\beta^0)} |\hat{\beta}_j - \beta_j^0| \\ &\leq (2\lambda(t) + \lambda) \|(\hat{\beta} - \beta^0)_{S_0}\|_1 . \end{aligned}$$

**Remark 5.** From Equation (15) we obtain the following inequality:

$$\begin{aligned} \|\hat{\beta}_{S_0^c}\|_1 &\leq \frac{\lambda + 2\lambda(t)}{\lambda - 2\lambda(t)} \|(\hat{\beta} - \beta^0)_{S_0^c}\|_1 \\ &\leq 3 \|(\hat{\beta} - \beta^0)_{S_0^c}\|_1 \end{aligned} \quad (16)$$

which allow us to use the compatibility condition.

Lemma 4.1, Condition 2.6 and the compatibility condition leads to:

$$\begin{aligned} (2\lambda(t) + \lambda) \|(\hat{\beta} - \beta^0)_{S_0}\|_1 &\geq C_1 \|f_{\hat{\beta}} - f_{\beta^0}\|_2^2 \\ &\geq C_2 C_1 \|x^T(\hat{\beta} - \beta^0)\|_2^2 \\ &\geq \frac{\phi_0^2}{C s_{\beta^0}} \|(\hat{\beta} - \beta^0)_{S_0}\|_1^2, \end{aligned}$$

where  $C := 1/(C_1 C_2)$ . Resuming we have

$$\frac{\phi_0^2}{s_{\beta^0} C} \|(\hat{\beta} - \beta^0)_{S_0}\|_1^2 \leq (2\lambda(t) + \lambda) \|(\hat{\beta} - \beta^0)_{S_0}\|_1.$$

(16) implies  $\|\hat{\beta} - \beta^0\|_1 \leq 4 \|(\hat{\beta} - \beta^0)_{S_0}\|_1$ , which yields

$$\frac{\phi_0^2}{4 s_{\beta^0} C} \|(\hat{\beta} - \beta^0)_{S_0}\|_1 \leq (2\lambda(t) + \lambda).$$

So we have

$$\|(\hat{\beta} - \beta^0)_{S_0}\|_1 \leq \lambda \frac{6 s_{\beta^0} C}{\phi_0^2}$$

and

$$\mathcal{E}(f_{\hat{\beta}}) \leq \lambda^2 \frac{9 s_{\beta^0} C}{\phi_0^2}.$$

□

## 5 Numerical results

In this section we present the result of a simulation study. We compare the following three methods:

- Our estimator, the censored regression with  $l_1$ -penalisation (CL).

$$\hat{\beta} = \arg \min_{\beta \in \mathcal{B}} \left\{ \frac{1}{n} \sum_{i=1}^n |y_i - \max\{c_i, x_i \beta\}| + \lambda \|\beta\|_1 \right\},$$

where  $\mathcal{B}$  is a large compact set, (See also (5).)

- The non-censored regression with  $l_1$ -penalty (NL).

$$\hat{\beta}^{NL} := \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n |y_i - x_i \beta| + \lambda \|\beta\|_1 \right\}.$$

This corresponds to taking the dataset and directly doing an  $l_1$ , non-censored regression with  $l_1$ -penalty on the entire dataset. That is, we consider the censored data as being non-censored data.

- Restricted dataset non-censored regression with  $l_1$ -penalty (RL).  
We take the non-censored data as a new restricted dataset and then do an  $l_1$ , non-censored regression with  $l_1$ -penalty on this restricted dataset. We have the following estimator:

$$\hat{\beta}^{RL} := \arg \min_{\beta} \left\{ \frac{1}{\#J} \sum_{i \in J} |y_i - x_i \beta| + \lambda \|\beta\|_1 \right\},$$

where  $J \subseteq \{1, \dots, n\}$  is the set of indexes corresponding to non-censored values of  $Y$  (i.e.  $i \in J \Leftrightarrow Y_i > c_i$ ).

### 5.1 Designs and settings

In our study we fit 24 different designs varying the following parameters

- $n$ , the number of observations. It is chosen to be 20, 30, 50 or 70.
- The dimension  $p$ . This is the number of parameters in our model.  $p$  can be 30, 50, 100 or 250.
- The sparsity  $s_{\beta^0}$ . This denotes the number of non-zero components of  $\beta^0$ . The sparsity  $s_{\beta^0}$  can be 3, 5 or 10.
- The signal to noise ratio (SNR). It is defined as

$$SNR := \sqrt{\frac{\sum_{i=1}^n (x_i \beta^0 \vee 0)^2}{\sum_{i=1}^n \varepsilon_i^2}}.$$

We take SNR to be 2 or 8.

The censoring factor  $c$  is random normally distributed with mean 0 and standard deviation 2.

We simulate the dataset as follows: We first randomly generate the  $n \times p$  dimensional matrix  $X$  (each component of the matrix is an independent realisation of a standard normal distributed random variable). Without loss of generality we take the active set of  $\beta^0$  as its first  $s_{\beta^0}$  components. These components take value  $\pm 1$  with probability 1/2 each. We take Gaussian 0-mean distributed errors, where the choice of the variance is implicitly given by the SNR. Finally the response variable is then generated by

$$y_i = \max\{x_i \beta^0 + \varepsilon_i, c_i\}$$

The tuning parameter  $\lambda$  is chosen to be  $\lambda = 8\sqrt{\log p/n}$ . This is (around) the theoretical choice given from Theorem 3.1.

In all our simulations about the same percentage of the data is censored. By construction we expect 50% of the data to be censored.

For the three different estimators we compare the prediction error

$$\frac{1}{n} \sum_{i=1}^n (x_i \hat{\beta} - x_i \beta^0)^2$$

and the estimation error

$$\|\hat{\beta} - \beta^0\|_1 .$$

The results of the simulation are summarised in Table 1

## 5.2 Conclusion

First of all we can notice that the NL method works much worse than the other methods in all possible designs for both the prediction and estimation error. Comparing CL with RL one can notice that our estimator works almost in any case (slightly) better than RL. This is not surprising because our estimator also takes into account the censored data.

**Remark 6** (High quantity of censored level). *As suggested in Condition 2.6, taking  $c$  constant equal 0 does not affect the quality of the fit. (Simulation results for this statement are not included in the paper because they just almost replicate Table 1).*

**Remark 7** (High quantity of censored data). *Increasing the number of censored components, for example changing the censoring level will clearly reduce the quality of all three estimators. This is not surprising because increasing the number of censored values somehow reduce the information we have in the observations. But as one can expect, our estimator is less affected than the other two estimators by such an increase. This also makes sense, because our estimator is the only one which is somehow able to deal with censored values.*

In conclusion one can summarize the results of the simulation study as follows:

- The simulation seem to confirm the theoretical results.
- Do not treat censored data as uncensored, it is better to exclude them from your dataset.
- Use an appropriate estimator (like CL) for analysing datasets containing censored data.
- The quality of the fit is strongly depending on the percentage of the censored data, but is almost not affected by the distribution of the censoring factor  $c$  (since the number of censored data does not change).

Simulation results							
$n^\circ$	setting	Estimation error			Prediction error		
$N$	$n, p, s_{\beta^0}, \text{STN}$	CL	NL	RL	CL	NL	RL
1	70,250,10,8	1.21 (0.5)	2.72 (0.34)	1.8 (0.55)	8.13 (3.12)	20.92 (1.49)	11.7 (2.94)
2	70,100,10,8	0.5 (0.18)	2.22 (0.33)	0.84 (0.56)	3.46 (1.39)	20.65 (2.49)	5.29 (3.18)
3	70,250,10,2	1.81 (0.22)	2.99 (0.32)	2.35 (0.31)	11.26 (1.62)	22.24 (1.63)	14.26 (1.43)
4	70,100,10,2	1.48 (0.2)	2.52 (0.31)	1.88 (0.32)	9.57 (1.93)	23.23 (2.35)	11.55 (2.21)
5	70,250,5,8	0.37 (0.19)	2.14 (0.29)	0.52 (0.48)	1.83 (0.92)	13.97 (1)	2.52 (2.27)
6	70,100,5,8	0.24 (0.04)	1.71 (0.22)	0.26 (0.06)	1.42 (0.25)	15.88 (2.36)	1.49 (0.32)
7	70,250,5,2	0.95 (0.15)	2.4 (0.24)	1.27 (0.3)	4.82 (0.94)	15.89 (1.38)	6.22 (1.52)
8	70,100,5,2	0.87 (0.12)	2.02 (0.22)	0.98 (0.21)	4.91 (0.75)	18.27 (2.14)	5.25 (1.13)
9	40,100,5,8	0.81 (0.44)	2.14 (0.28)	1.04 (0.56)	3.74 (2.23)	12.41 (1.25)	4.5 (2.35)
10	40,100,5,2	1.2 (0.21)	2.37 (0.32)	1.57 (0.42)	5.35 (1.17)	14.09 (1.78)	6.44 (1.52)
11	40,100,3,8	0.25 (0.15)	1.38 (0.22)	0.28 (0.27)	1.15 (0.82)	7.79 (1.21)	1.26 (1.28)
12	40,100,3,2	0.7 (0.14)	1.65 (0.25)	0.76 (0.2)	3.06 (0.81)	9.13 (1.17)	3.13 (0.94)
13	40,50,5,8	0.5 (0.35)	2.07 (0.39)	0.47 (0.33)	2.13 (1.7)	11.2 (1.77)	1.97 (1.51)
14	40,50,5,2	1.07 (0.27)	2.01 (0.35)	1.24 (0.36)	4.68 (1.12)	11.25 (2.06)	5.37 (1.39)
15	40,50,3,8	0.2 (0.11)	1.62 (0.26)	0.2 (0.11)	0.7 (0.36)	8.17 (1.32)	0.73 (0.4)
16	40,50,3,2	0.75 (0.18)	1.75 (0.28)	0.83 (0.21)	2.75 (0.75)	8.86 (1.41)	3.03 (0.76)
17	40,50,5,8	0.45 (0.27)	1.91 (0.35)	0.62 (0.44)	1.93 (1.35)	10.45 (1.72)	2.62 (1.93)
18	40,50,5,2	1.1 (0.26)	2.21 (0.35)	1.32 (0.41)	4.74 (1.48)	12.5 (1.99)	5.42 (1.88)
19	40,50,3,8	0.21 (0.11)	1.83 (0.37)	0.19 (0.07)	0.75 (0.36)	9.56 (1.88)	0.69 (0.27)
20	40,50,3,2	0.66 (0.2)	1.6 (0.29)	0.74 (0.22)	2.5 (0.81)	8.47 (1.2)	2.93 (0.92)
21	20,30,5,8	1.5 (0.33)	1.98 (0.45)	1.37 (0.49)	4.85 (0.99)	7.07 (1.37)	4.37 (1.24)
22	20,30,5,2	1.4 (0.33)	1.81 (0.45)	1.53 (0.42)	4.76 (0.96)	6.48 (1.03)	4.88 (0.97)
23	20,30,3,8	0.58 (0.4)	1.28 (0.35)	0.67 (0.5)	1.75 (1.27)	4.36 (1.25)	1.96 (1.41)
24	20,30,3,2	1.26 (0.3)	1.7 (0.35)	1.16 (0.42)	2.99 (0.78)	5.09 (1.09)	2.93 (1.1)

Table 1: Results of the simulation study. For the 24 different designs ( $N = 1, \dots, 24$ ) the performance of our estimator (CL) is compared with (NL) and (RL). Based on 30 replicates per design the average and (in brackets) the standard deviation of the prediction and the estimation errors are given in the table.



## References

- Bickel, P. J., Ritov, Y., and Tsybakov, A. B. (2009). Simultaneous analysis of lasso and dantzig selector. *Annals of Statistics*, 37(4):1705–1732.
- Bühlmann, P. and van de Geer, S. (2011). *Statistics for High-Dimensional Data*. Springer.
- Chay, K. Y. and Powell, J. L. (2001). Semiparametric censored regression models. *The Journal of Economic Perspectives*, 15(4):pp. 29–42.
- Koltchinskii, V. (2009). The Dantzig selector and sparsity oracle inequalities. *Bernoulli*, 15:799–828.
- Koltchinskii, V. (2011). *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems: École d’Été de Probabilités de Saint-Flour XXXVIII-2008*, volume 2033. Springer.
- Murphy, S. A., van der Vaart, A. W., and Wellner, J. A. (1999). Current status regression. *Mathematical Methods of Statistics*, 8(3).
- Powell, J. L. (1984). Least absolute deviations estimation for the censored regression model. *Journal of Econometrics*, 25(3):303 – 325.
- Städler, N., Bühlmann, P., and van de Geer, S. (2010).  $l_1$ -penalization for mixture regression models. *Test*, page 246.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Ser. B*(58):267–288.
- van de Geer, S. (2000). *Empirical Processes in M-Estimation*. Cambridge University Press.
- van de Geer, S. (2007). The deterministic Lasso. In *JSM proceedings, 2007, 140*. American Statistical Association.
- Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B*, 67:301–320.